



# Машинное обучение для решения исследовательских и инженерных задач в науках о Земле

Михаил Криницкий

К.Т.Н., Н.С.

Институт океанологии РАН им. П.П. Ширшова

Лаборатория взаимодействия океана и атмосферы и  
мониторинга климатических изменений (ЛВОАМКИ)



# Классификация задач и методов машинного обучения

Михаил Криницкий

К.Т.Н., Н.С.

Институт океанологии РАН им. П.П. Ширшова

Лаборатория взаимодействия океана и атмосферы и  
мониторинга климатических изменений (ЛВОАМКИ)

# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

ЦЕЛЬ: сформулировать задачу (в терминах машинного обучения)

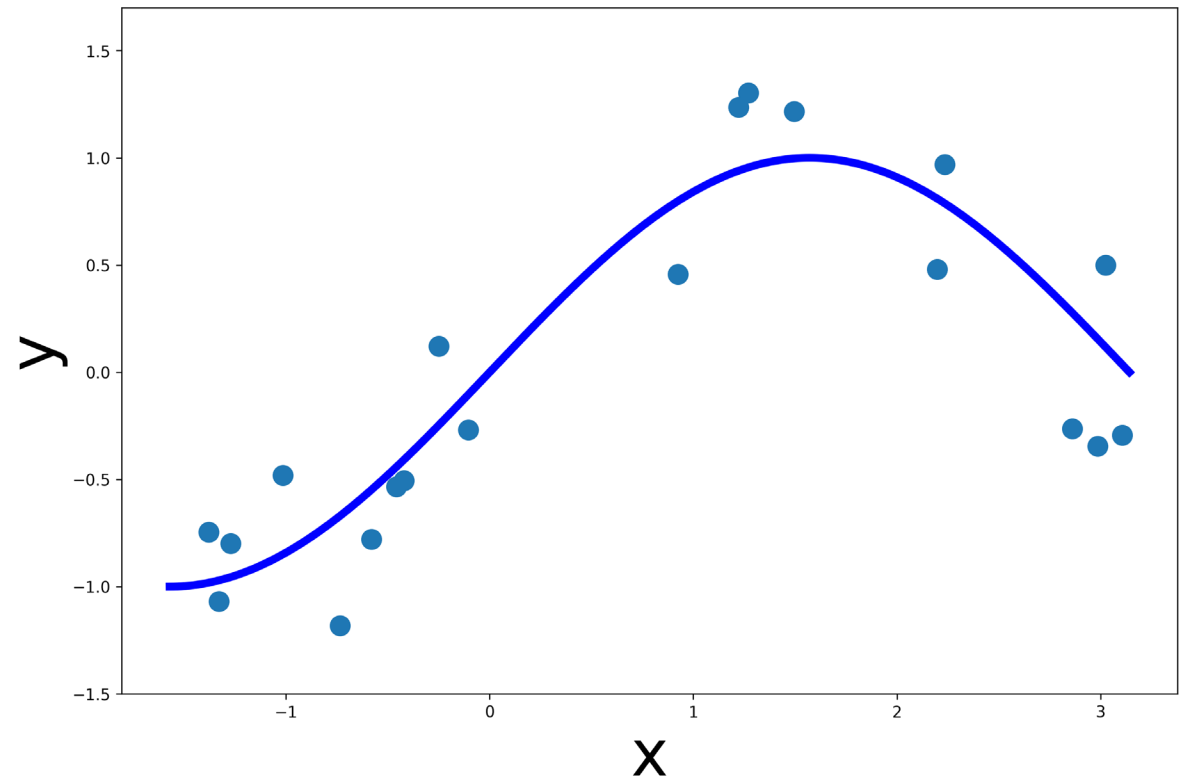
○ «Обучение с учителем»

- восстановление регрессии

**что я хочу?** – значение  $y$

$$y \in \mathbb{R}^m$$

$m$  – размерность целевой переменной



## ремарки

- Количество размеченных данных (зачастую) играет роль
- Разные модели ведут себя по-разному в зависимости от шума в данных, от количества данных, от наличия выбросов в данных
- Сложная точная модель – не обязательно лучшая для конкретной задачи

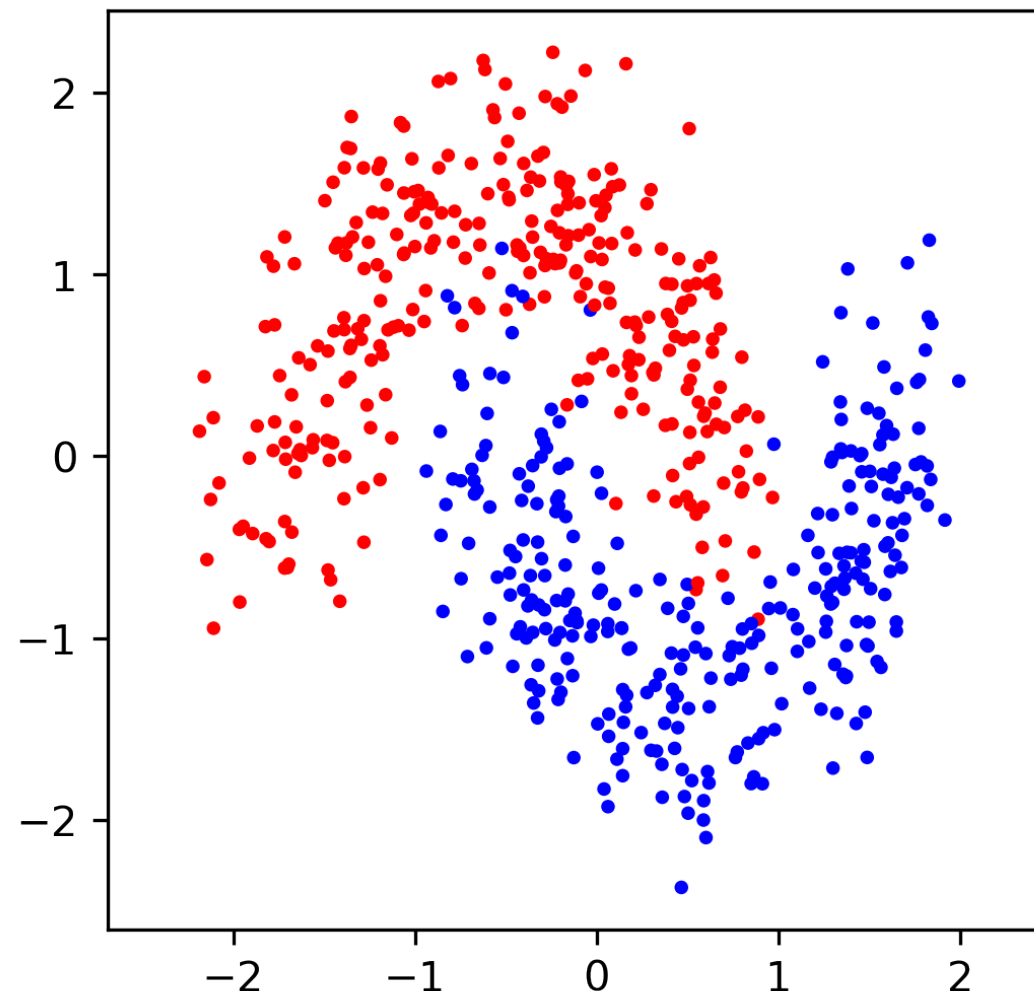
# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

ЦЕЛЬ: сформулировать задачу (в терминах машинного обучения)

○ «Обучение с учителем»

- восстановление регрессии
- классификация

**что я хочу?** — метку класса  
**«красный или синий?»**  
(бинарная классификация)



# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

ЦЕЛЬ: сформулировать задачу (в терминах машинного обучения)

типы задач:

○ «Обучение с учителем»

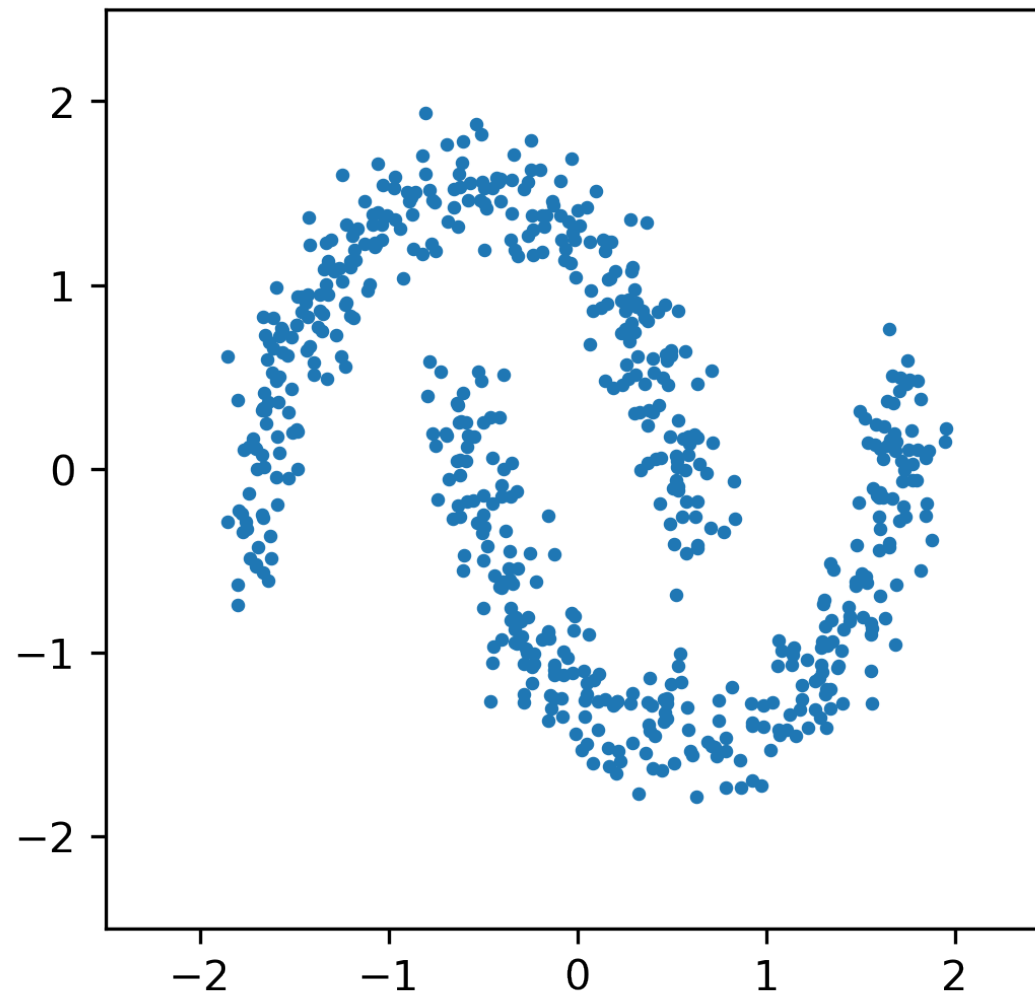
- восстановление регрессии
- классификация

○ «Обучение без учителя»

- поиск структуры в данных

**что я хочу?**

- метки групп
- знать, есть ли группы?
- сколько групп?



## ремарки

Кластеризация:

- Количество данных часто, но не всегда играет роль
- Разные модели ведут себя по-разному в зависимости от шума в данных, от количества данных, от наличия выбросов в данных, от наличия структуры в данных
- Разные модели дают разный результат, но нет «более правильного» результата. Есть «более подходящий» для целей конкретного исследования.

# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

ЦЕЛЬ: сформулировать задачу (в терминах машинного обучения)

типы задач:

## ○ «Обучение с учителем»

- восстановление регрессии
- классификация

## ○ «Обучение без учителя»

- кластеризация
- понижение размерности

## что я хочу?

- новое представление (признаковое описание) данных в пространстве меньшей размерности
- цели:
- Визуализация на плоскости, в 3D
  - Борьба с переобучением (в контексте т.н. «проклятия размерности»)
  - Сжатие данных с минимальными потерями
  - Сокращение вычислительных затрат при обработке данных
  - Извлечение значимых признаков, feature engineering

Пожелания:

- Сохранение структуры данных
- Сохранение отношений близости между объектами (событиями)
- Возможность визуализации
- Интерпретируемость новых признаков

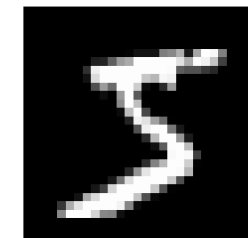
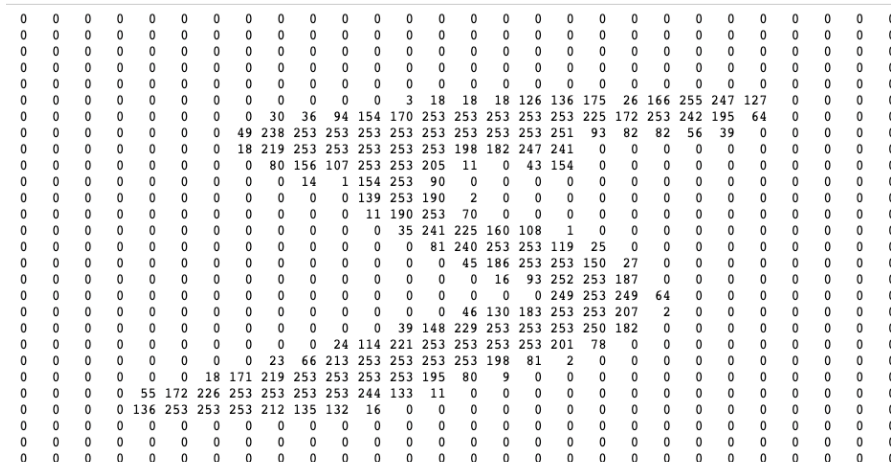
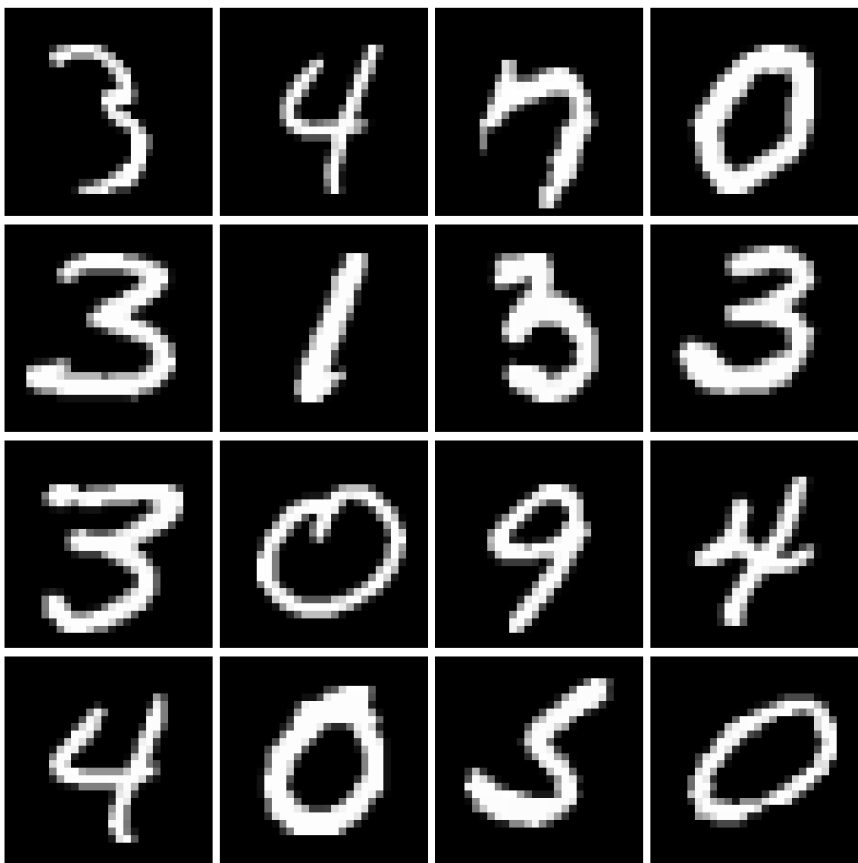


# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

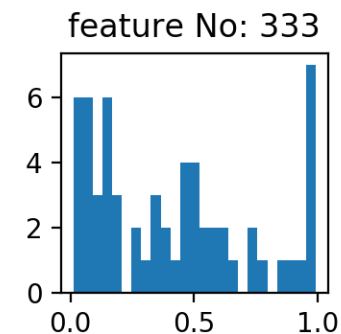
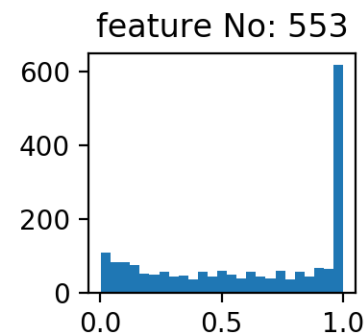
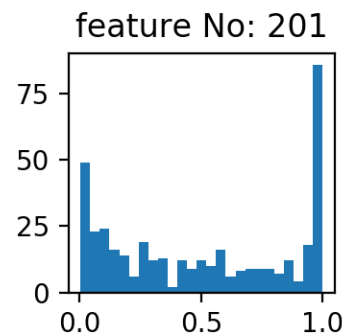
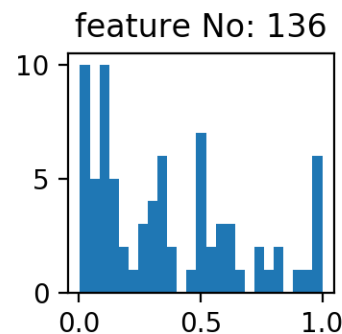
## Задача понижения размерности: пример

# MNIST dataset\*

example #0 (label = "5")



data distribution (4 of 784 features)



\* <http://yann.lecun.com/exdb/mnist/>

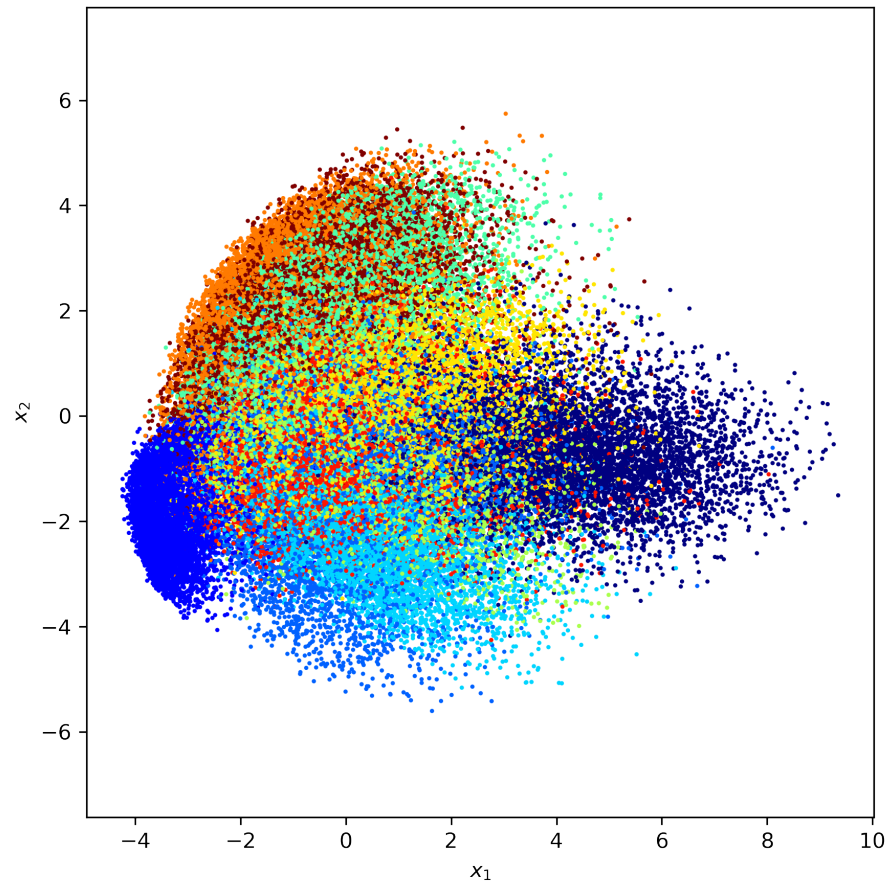
# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

## Задача понижения размерности: пример

### PCA

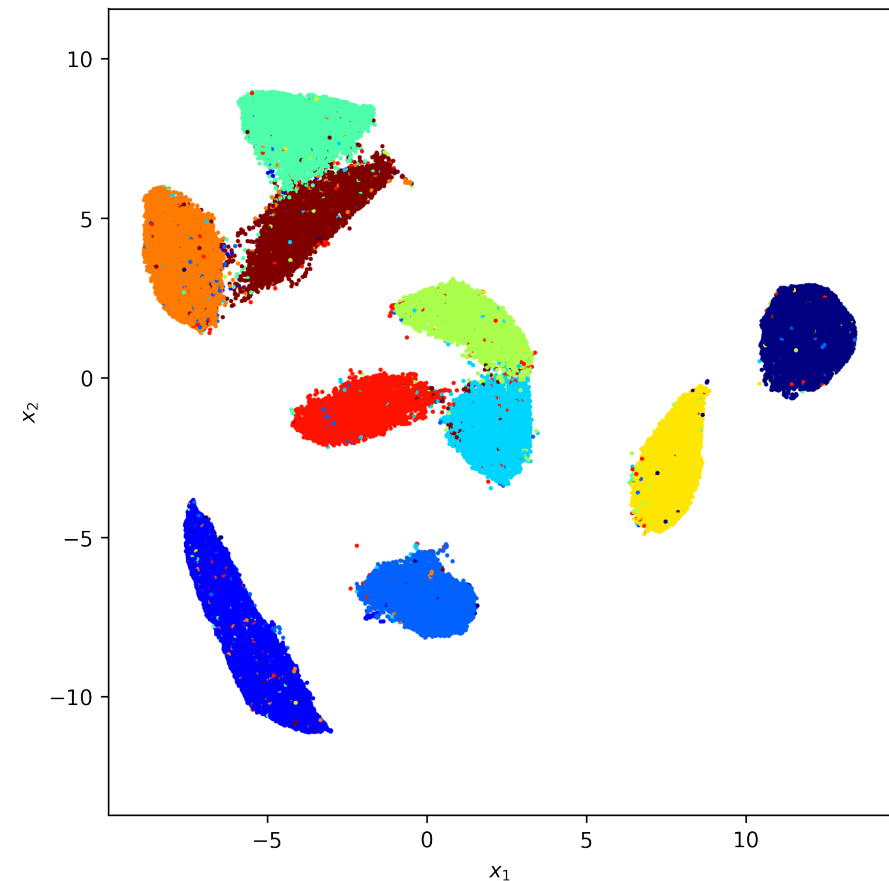
Principal components analysis

Метод главных компонент



### UMAP

Uniform Manifold Approximation and  
Projection



# ремарки

Понижение размерности:

- Количество данных почти всегда играет роль
- Разные модели ведут себя по-разному в зависимости от шума в данных, от количества данных, от наличия выбросов в данных, от наличия структуры в данных
- Разные модели дают разный результат, нет «более правильного» результата. Но есть «более подходящий» для целей конкретного исследования.
- Модели различаются по интерпретируемости (e.g. PCA vs. UMAP)

# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

типы задач:

- «Обучение с учителем»
  - восстановление регрессии
  - классификация
- «Обучение без учителя»
  - кластеризация
  - понижение размерности
  - восстановление распределения данных

**что я хочу?**

- Получить модель, генерирующую примеры, распределение которых совпадает с распределением обучающих данных

Цели:

- дополнение данных
- заполнение пропусков в данных

**Примеры:**

- DeepFake - video
- SuperResolution - images
- Text-to-speech - audio (Siri, Алиса, Cortana, Alexa etc.)

# КЛАССИФИКАЦИЯ ЗАДАЧ МАШИННОГО ОБУЧЕНИЯ

типы задач:

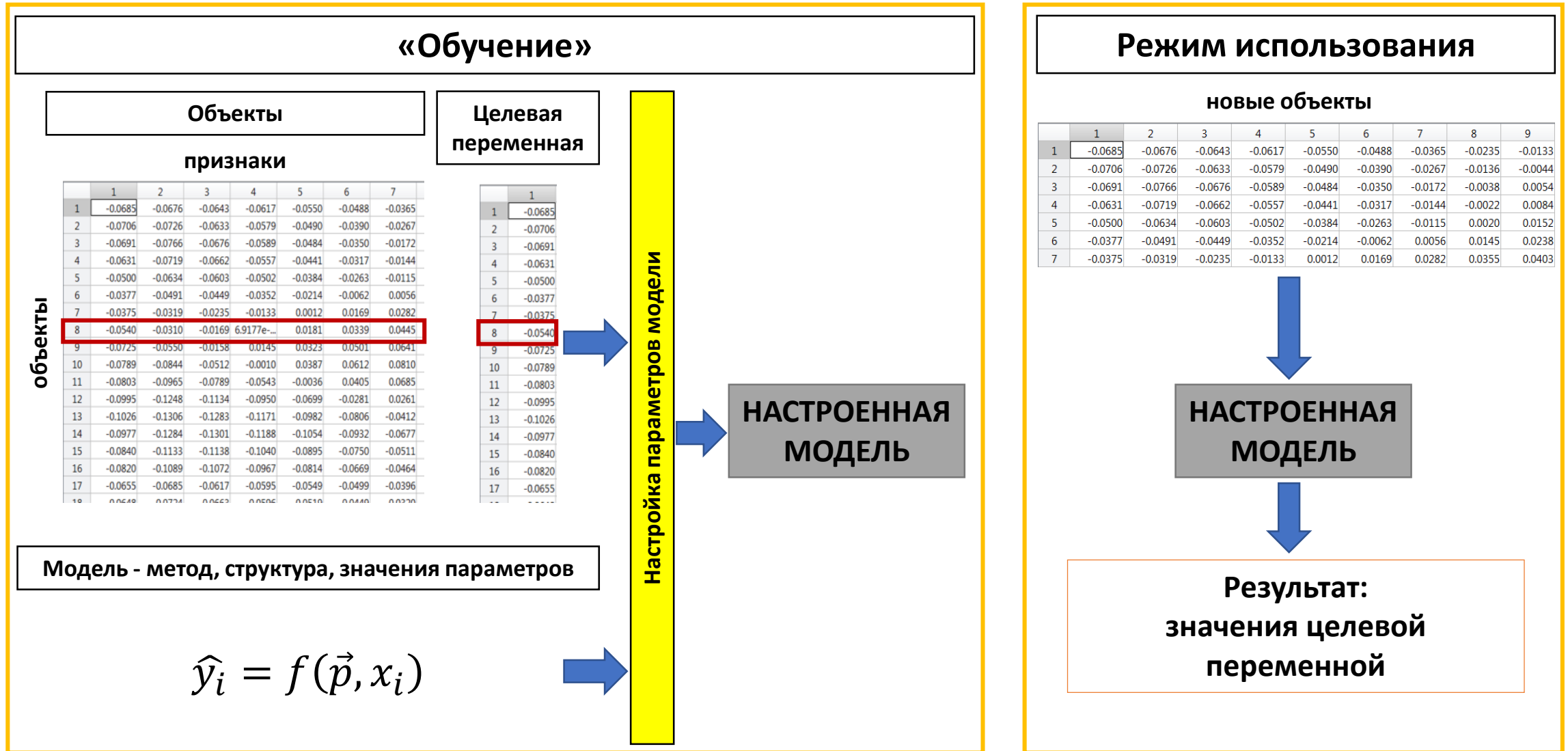
- «Обучение с учителем» - **Supervised learning**
  - восстановление регрессии
  - классификация
- «Обучение без учителя» - **Unsupervised learning**
  - кластеризация
  - понижение размерности
  - восстановление распределения данных

**другие (реже используемые в ES) типы**

- «Обучение с частичным привлечением учителя»  
- **Weakly supervised learning**
- «Обучение с подкреплением» -  
**Reinforcement learning**

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

обучаем (тренируем) модель на имеющихся данных



# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

обучаем (тренируем) модель на имеющихся данных

(вернее, на большей их части)

формулировка задачи:

$x \in X$  — объекты

$y \in Y$  — ответы

$\mathcal{F}: X \rightarrow Y$  — искомая закономерность

$\mathcal{T}: \{x_i; y_i\}$  — «обучающая выборка»  
(прецеденты)

Найти:  $\hat{\mathcal{F}}: \{x_i\} \rightarrow \{y_i\}$

# ОБЩАЯ СХЕМА РЕШЕНИЯ ЗАДАЧ ОБУЧЕНИЯ С УЧИТЕЛЕМ

обучаем (тренируем) модель на имеющихся данных

формулировка задачи:

$x \in \mathbb{X}$  — объекты

$y \in \mathbb{Y}$  — ответы

$\mathcal{F}: \mathbb{X} \rightarrow \mathbb{Y}$  — искомая закономерность

$\mathcal{T}: \{x_i; y_i\}$  — «обучающая выборка»  
(прецеденты)

Найти:  $\hat{\mathcal{F}}: \{x_i\} \rightarrow \{y_i\}$

один из способов решения:

$\mathcal{L}(\hat{\mathcal{F}}(x))$  — функционал ошибки  
(эмпирического риска, потерь)

$\hat{y}_i = \hat{\mathcal{F}}(x_i) = f(\vec{p}, x_i)$  — функционально задаваемая зависимость. **Предположение исследователя о виде закономерности.** Иногда задается параметрически,  $\vec{p}$  — вектор параметров.

$\mathcal{L} = L(\vec{p}, \mathcal{T})$  — функция ошибки

$$\hat{p} = \underset{\mathbb{P}}{\operatorname{argmin}}(L(\vec{p}, \mathcal{T}))$$

$$\hat{\mathcal{F}} = f(\hat{p}, x)$$



# ЛИНЕЙНАЯ РЕГРЕССИЯ